# KIOXIA

# RAID Offload and Its Application
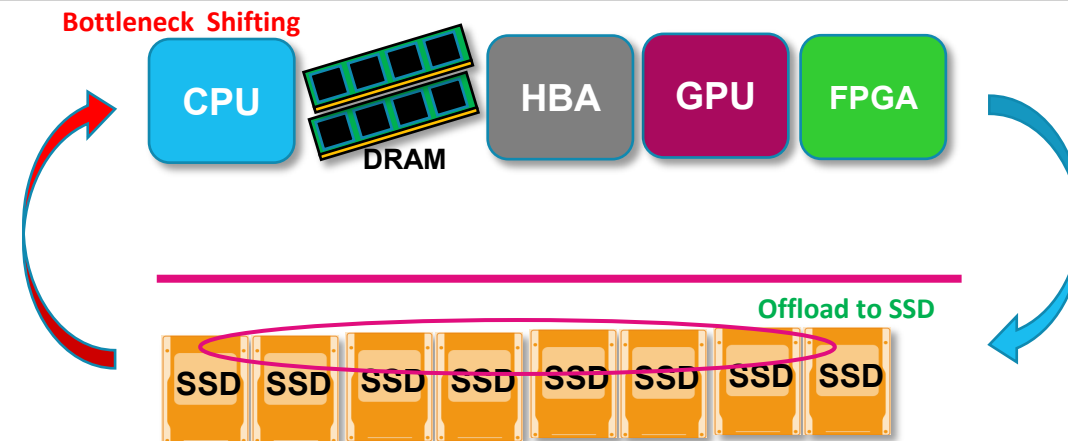
**Aug 6, 2024**

**By Chandra Nelogal, Dell®**

**Devesh Rai, KIOXIA**

DELL Technologies

# RAID Offload Introduction

- NVMe™ SSDs performance improvement is continuously shifting the bottlenecks to applications

- Industry is addressing the problem in different ways

  - Hardware - RAID host bus adapters (HBA), data processing units (DPU), field programmable gate array (FPGA), configurable spatial accelerator (CSA)

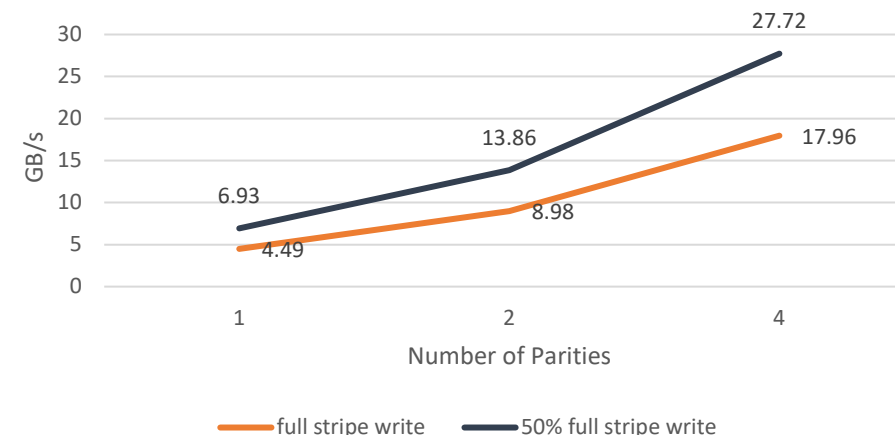  - Software - Costly CPU cores and dedicated DRAM

- **KIOXIA proposes to offload:**

  - Parity compute and memory resources to SSD

- **Benefits**

  - Reducing system total cost of acquisition (TCA) and total cost of ownership (TCO)

  - Continue to leveraging existing RAID applications and fault management



**Bottleneck Shifting**

CPU  DRAM  HBA  GPU  FPGA

Offload to SSD

SSD SSD SSD SSD SSD SSD SSD SSD

**DRAM Throughput for RAID Write @ 1 Gigabytes per Second (GB/s)**



- full stripe write: 4.49, 8.98, 17.96
- 50% full stripe write: 6.93, 13.86, 27.72

Number of Parities: 1, 2, 4
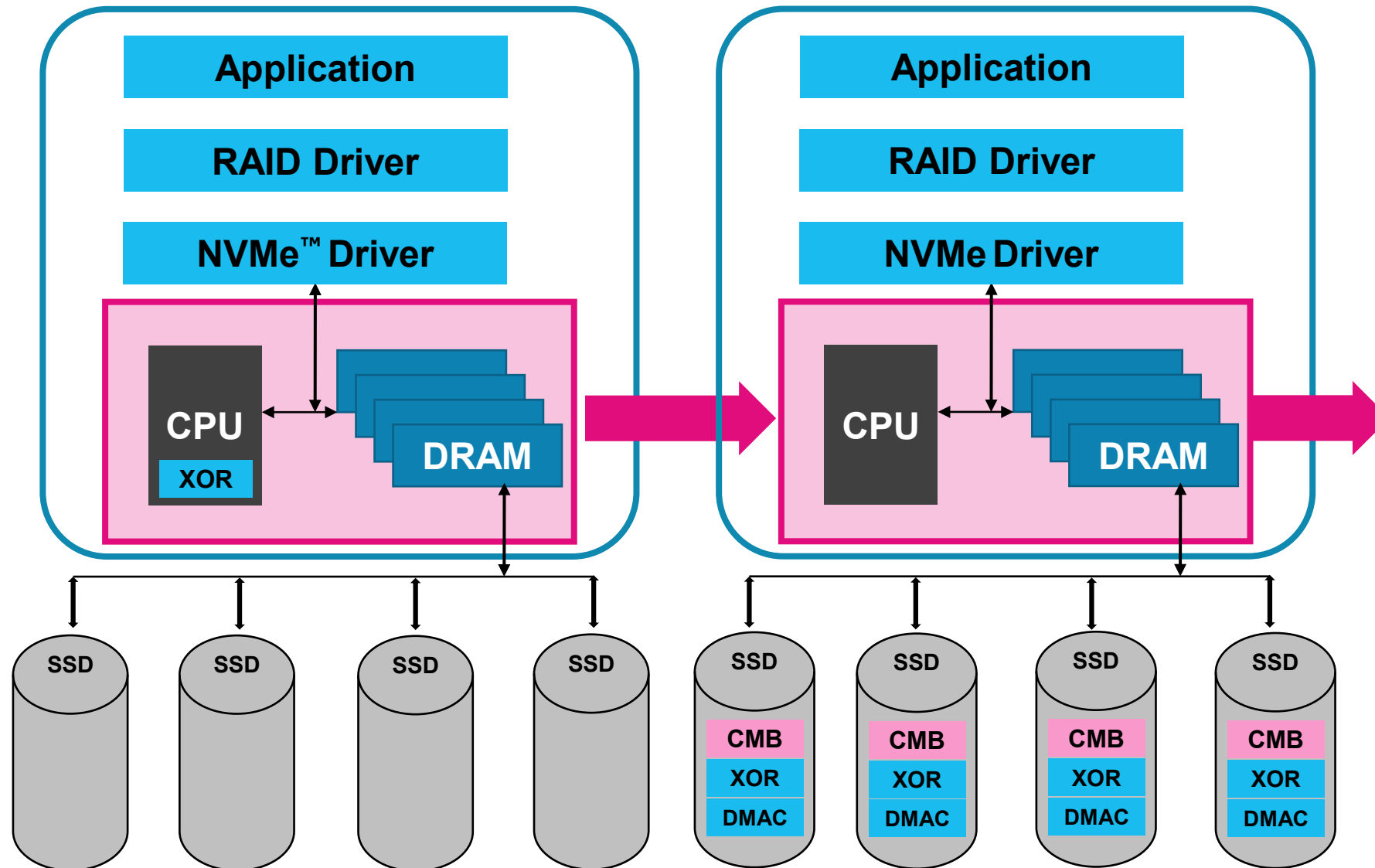
— full stripe write   — 50% full stripe write

\* Implementation specific

Images and product icons created by KIOXIA
Chart source:  KIOXIA strategic marketing in-house testing and calculation
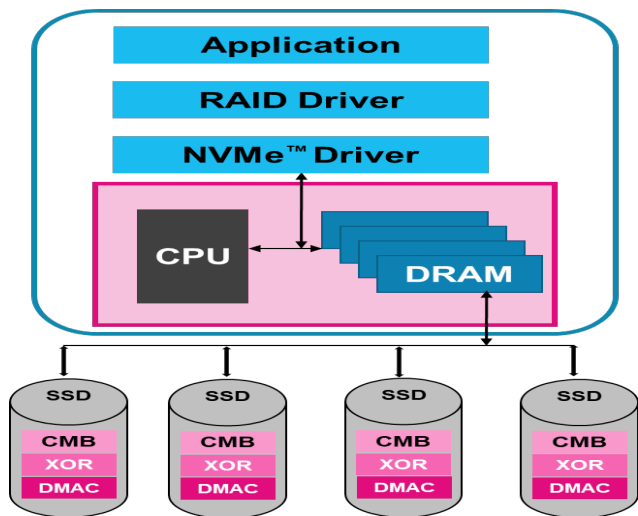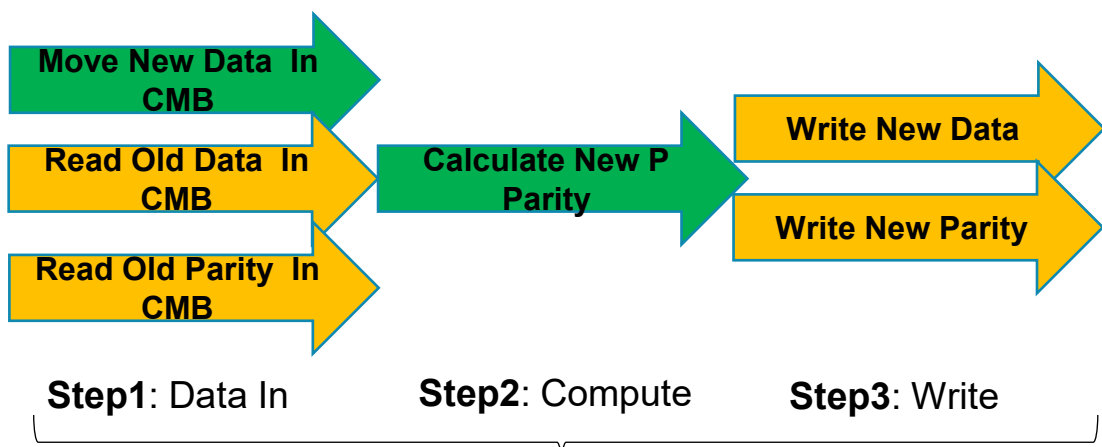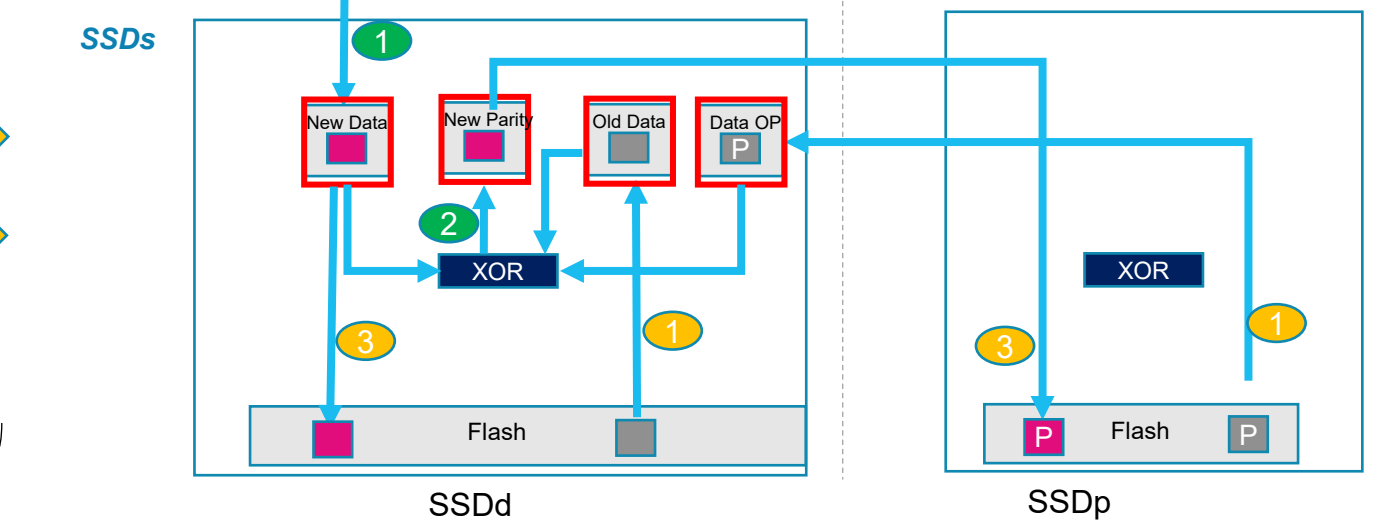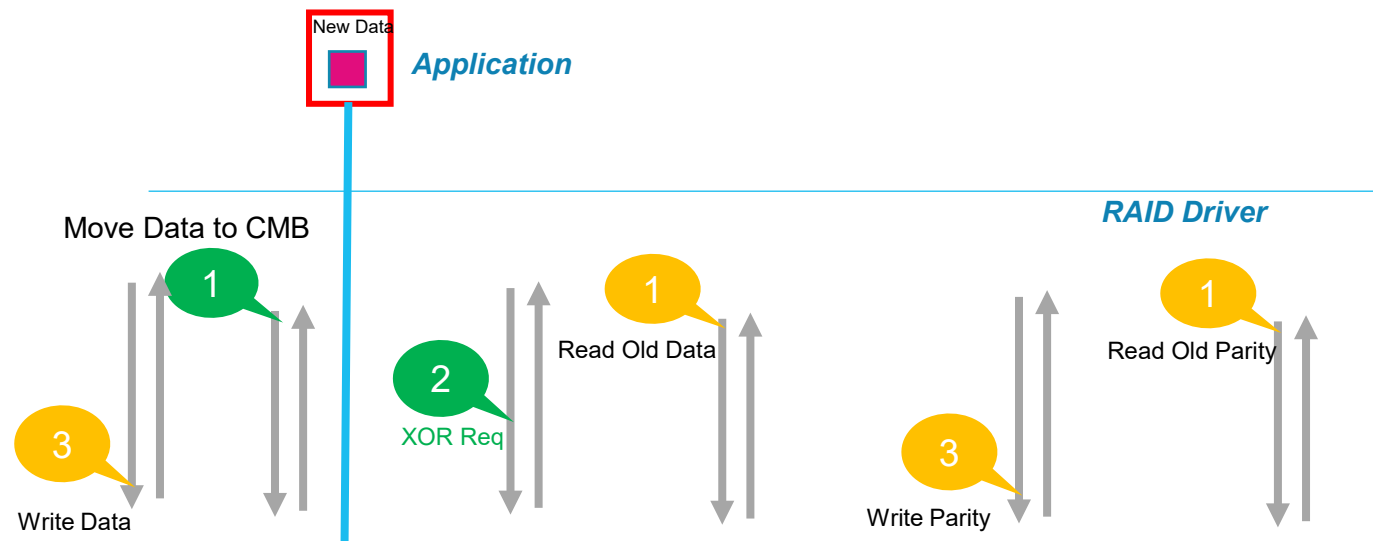
# KIOXIA RAID Offload Resource Utilization



- Host managed standard based offload solution

- KIOXIA NVMe™ SSDs feature:
  - Controller memory buffers (CMB) – for DRAM offload
  - Exclusive OR (XOR) – up to 8 parity compute
  - Direct memory access controller (DMAC) – to place data in host address space (including remote CMB)

- Parallel compute and linear scaling

Graphics and product icons created by KIOXIA

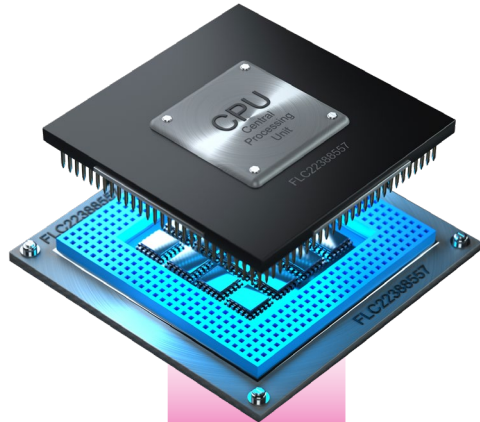# Command and Data Flow Example for RAID 5 Write

# Proof of Concept (PoC) with mdraid5 and KIOXIA CM7 Series SSD

1 Reduce CPU workload for RAID computation

2 Reduce DRAM bandwidth Utilization

3 Improve host CPU utilization; contribute to energy efficiency for PCIe® Gen 5 server

**Leverage High Performance of PCIe® 5.0 SoCs**

RAID 5

Parity Calculation Offload

## RAID Offload : PoC Results (with KIOXIA CM7 & mdraid5)

| System | KIOXIA CM7 Gen4 x4 – mdRAID 5# | RAID Offload | % Benefit |
|---|---|---|---|
| Number of SSDs | 5 | 5 | |
| Full Stripe Write 512 kibibytes (KiB) | | | |
| CPU Utilization | 42 | 37 | 12% Reduction |
| DRAM Bandwidth in mebibytes (MiB/s) | 3450 | 340 | 91% Reduction |

IO workload: Flexible I/O tester (FIO) 512K Random Write @ 950 megabytes per second (MB/s)

System DELL® PowerEdge™ R650xs Xeon™ Gold 6338N 2.2GHz(2 Socket, 32 Cores) PCIe Gen4 , SSDs : 5xCM7 Gen4 (1.92TB)
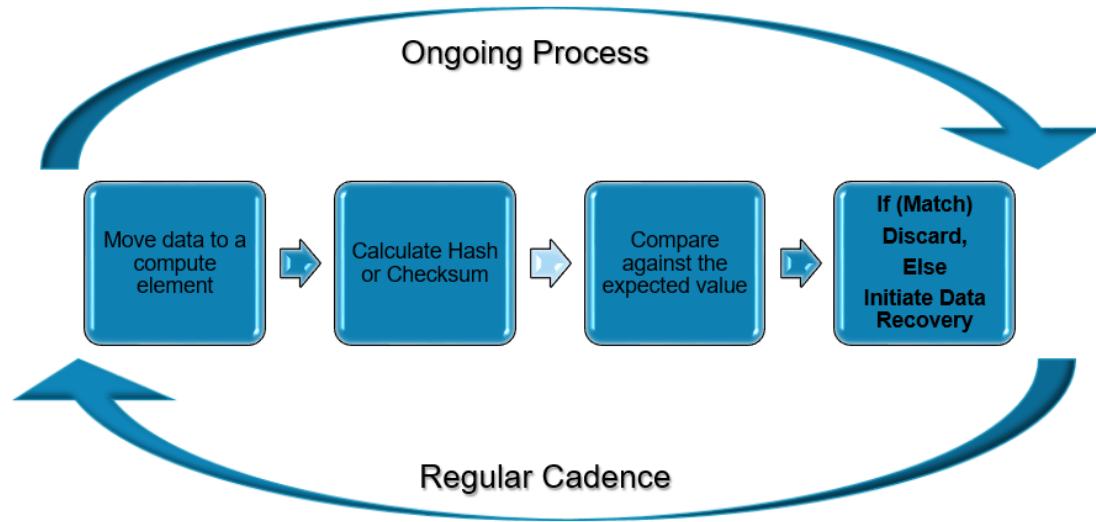
The KIOXIA product images shown are a representation of the design model and not an accurate product depiction.

KIOXIA
CM7
BiCS FLASH™
Enterprise NVMe™ SSD

Image source:

# RAID Offload Use Cases

# Data Scrubbing in Conventional Setup

- **Data Scrubbing: early detection and correction of errors**

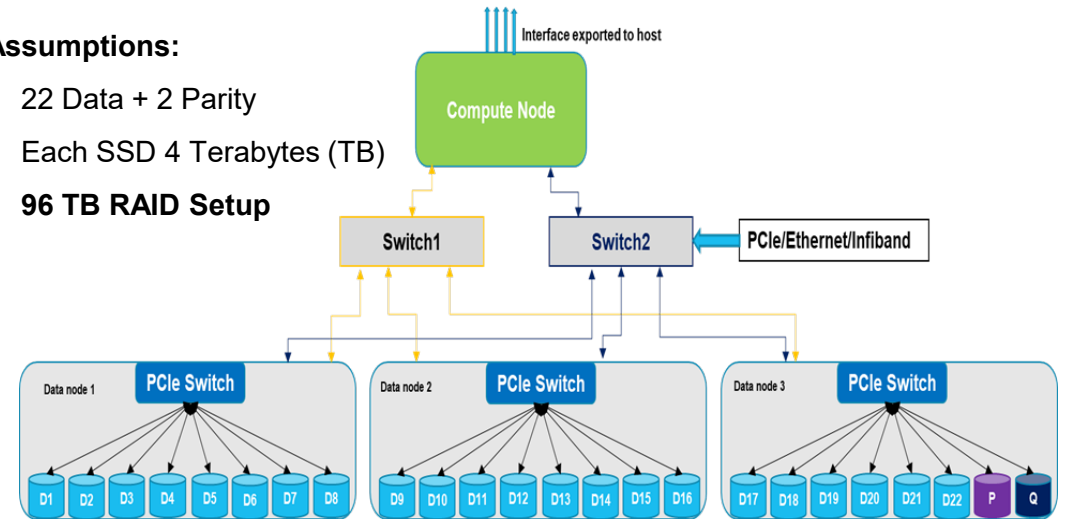- **Data Scrubbing technology: hash, checksum or RAID technology**



Ongoing Process

Move data to a compute element → Calculate Hash or Checksum → Compare against the expected value → If (Match) Discard, Else Initiate Data Recovery

Regular Cadence

- **All data movement during scrubbing operation is an overhead penalty paid to ensure data integrity**

**Assumptions:**

- 22 Data + 2 Parity

- Each SSD 4 Terabytes (TB)

- **96 TB RAID Setup**



Interface exported to host

Compute Node

Switch1    Switch2    PCIe/Ethernet/Infiband

Data node 1  PCIe Switch    Data node 2  PCIe Switch    Data node 3  PCIe Switch

D1 D2 D3 D4 D5 D6 D7 D8    D9 D10 D11 D12 D13 D14 D15 D16    D17 D18 D19 D20 D21 D22 P Q

**Compute node performing disk scrubbing for one stripe using RAID**

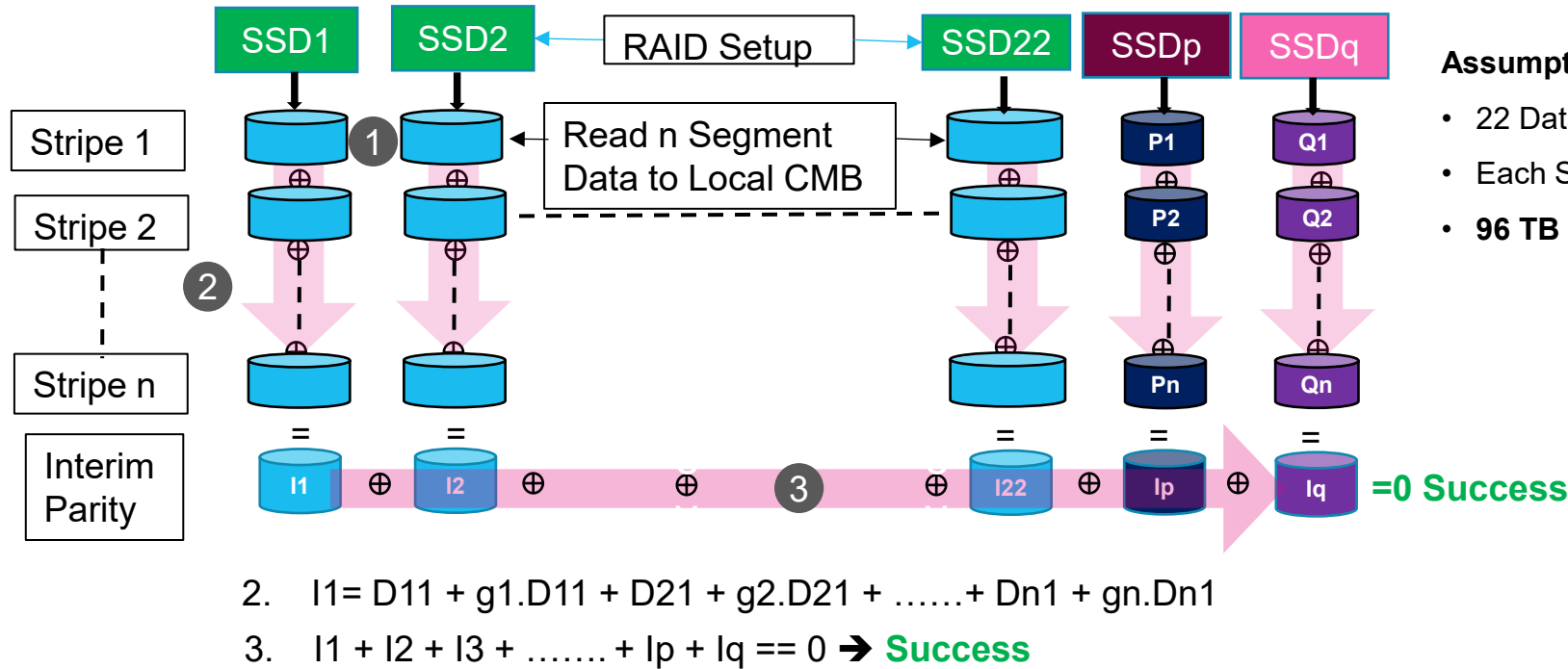$P + D1 + D2 + D3 + D4 \ldots\ldots + D22 = 0$

$Q + g1.D1 + g2.D2 + g3.D3 + \ldots.. + g22.D22 = 0$

**In above setup, 96TB data moves over PCIe®, network, CPU and 192TB through memory subsystem during each scrubbing cycle**

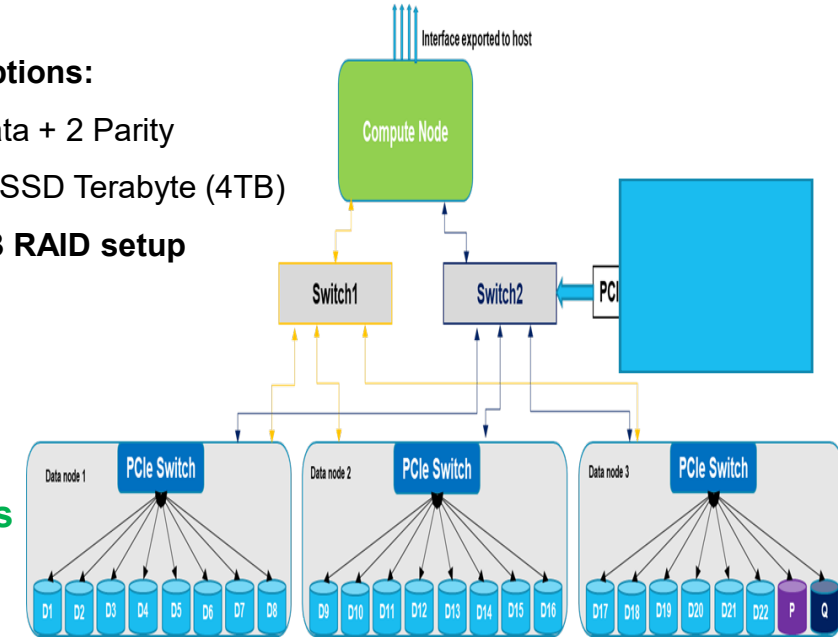Graphics and product icons created by KIOXIA

Assumptions created by KIOXIA in-house engineering team

# Data Scrubbing using RAID/EC Offload

**Assumptions:**

- 22 Data + 2 Parity
- Each SSD Terabyte (4TB)
- **96 TB RAID setup**

2. $I1 = D11 + g1.D11 + D21 + g2.D21 + \ldots\ldots + Dn1 + gn.Dn1$

3. $I1 + I2 + I3 + \ldots\ldots + Ip + Iq == 0$ ➔ **Success**
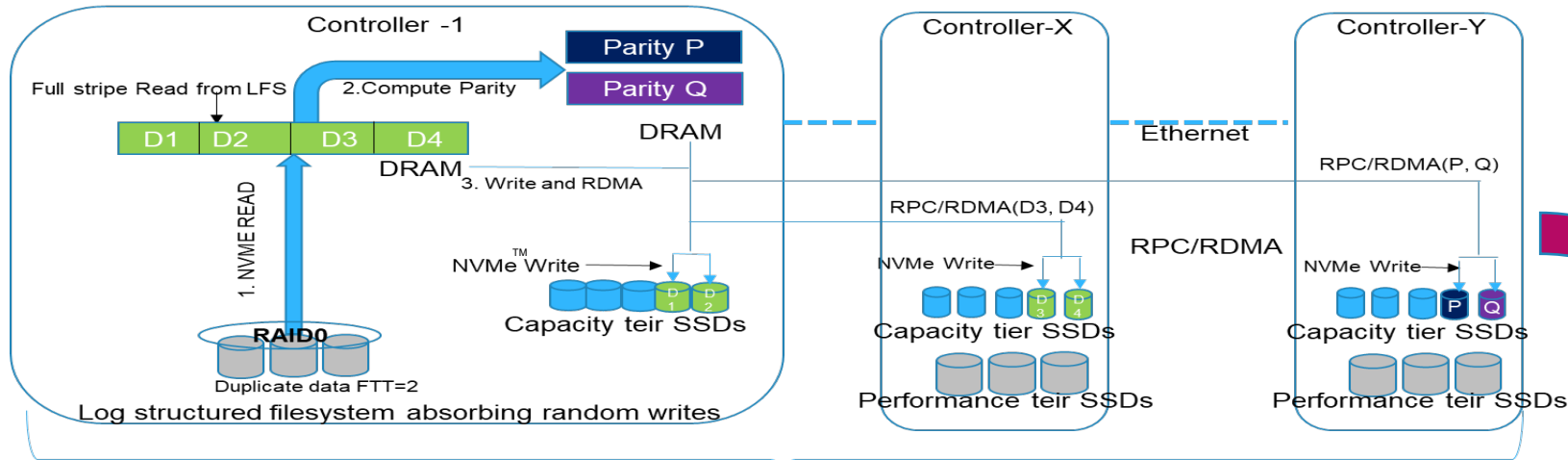
❖ Using 3 step process, ~99% data movement can be reduced

❖ No data passes through CPU and DRAM on compute node

❖ For n stripes, only one stripe moves over network and PCIe[®]

❖ Data scrubbing proof of concept data shown in table is for 9 SSD

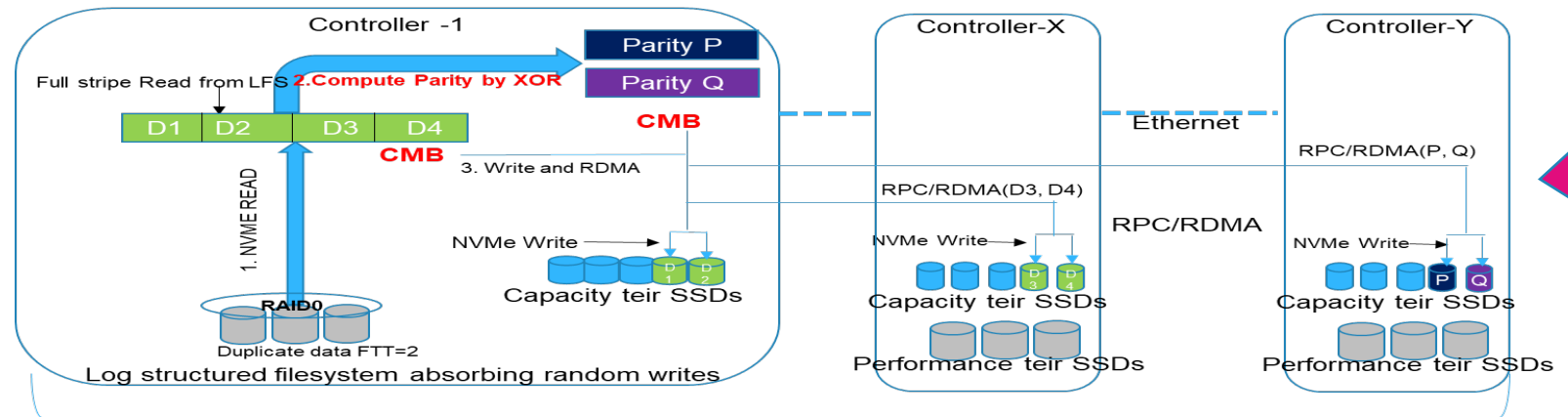| Resource Utilization | Offload Disabled | Offload Enabled |
|---|---|---|
| Scrubbing time | 129s | 91s |
| DRAM Bandwidth | 10.24 GB/s | 1.43 GB/s |
| Total CPU Utilization | 99.5% | ~70% |
| L3 Cache Misses | 14.7M | 4M |
| Total PCIe Write (MB/s) | 3694 MB/s | 159 MB/s |

Graphics, product icons, and tables created by KIOXIA. Assumptions created by KIOXIA in-house engineering team

# RAID Offload in Hyperconverged, Software Defined Storage (SDS)

# Future Possibilities: A Call to Action

- Offloading data scrubbing on to SSDs can significantly alleviate memory and network bandwidth bottlenecks and reduce data movement

  - Better resource utilization

- Easily adoptable in existing hyper converged infrastructure (HCI), RAID or similar solutions

- Dell$^®$ is collaborating with KIOXIA to standardize this technology

**Standards Based**

**Host Controlled**

**Hardware Accelerators**
(Memory, Compute, DMAC)

## Let's Collaborate! Visit Booth# 307 for Demo

# Example : Parity P and Q Generation Mechanism

RAID 6 P Parity:      $P = D0 \oplus D1 \oplus D2 \ldots \oplus D31$

RAID 6 Q Parity:      $Q = g0 \cdot D0 \oplus g1 \cdot D1 \oplus g2 \cdot D2 \ldots \oplus g31 \cdot D31$

1. g0 … g31 are Galois coefficient provided by host in XOR command.

2. D0 … D31 are per SSD data segment in a given full stripe

3. RAID application has option to calculate P or calculate Q or calculate both

4. Proposing up to 8 parity request in single command to support erasure code

5. Command structure may change during standardization process

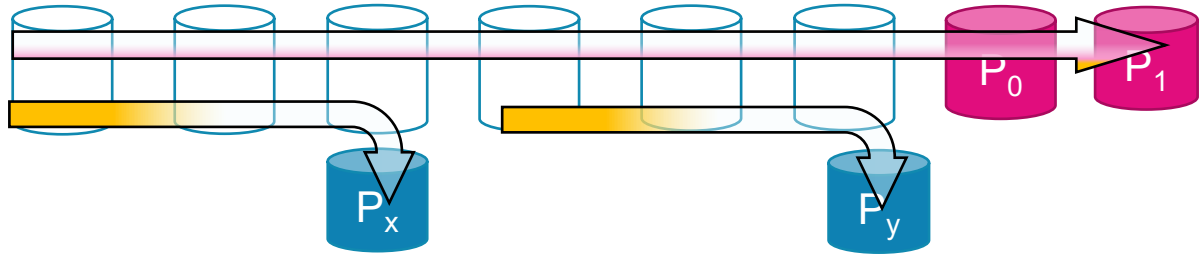| Field for P parity command | Value |
|---|---|
| Source buffer address | D0, D1, ……D31 |
| Galois coefficient for each buffer | 1 |
| Each source buffer length | 16 Kilobytes (KB) |
| Output buffer address | P |
| Number of source buffer | 32 |

| Field for Q parity command | Value |
|---|---|
| Source buffer address | D0, D1, ……D31 |
| Galois coefficient for each buffer | g0, g1, ……..g31 |
| Each source buffer length | 16 Kilobytes (KB) |
| Output buffer address | Q |
| Number of source buffer | 32 |

Single XOR command calculating P and Q parity

Tables: created by KIOXIA

KIOXIA

Drives

- **Erasure Code** command for **4 parity** compute

| Parity $P_x$ | |
|---|---|
| Src buf | x0, x1, x2 |
| Galois coefficient | 1,1,1,1 |
| Output buffer address | Px |
| Operation type | XOR |

| Parity $P_y$ | |
|---|---|
| Src buf | y0,y1,y2 |
| Galois generator | 1,1,1,1 |
| Output buffer address | Py |
| Operation type | XOR |

| Parity $P_0$ | |
|---|---|
| Src buf | x0, x1, x2,y0,y1,y2 |
| Galois coeffficient | $\alpha_0, \alpha_1, \alpha_2, \beta_0, \beta_1, \beta_2$ |
| Output buffer address | $P_0$ |
| Operation type | XOR |

| Parity $P_1$ | |
|---|---|
| Src buf | x0, x1, x2,y0,y1,y2 |
| Galois coefficient | $\alpha_0^2, \alpha_1^2, \alpha_2^2, \beta_0^2, \beta_1^2, \beta_2^2$ |
| Output buffer address | $P_1$ |
| Operation type | XOR |